

1. Figure 1.4 of your text shows the basics while Figure 2.21 shows abt more detail. This HW is over chapter 1 so let's talk about the basics.

The compiler translates the high-level programming language such as C/C++ into ^{an} assembly language program specific to the target ISA, such as MIPS.

The assembler takes the assembly language program and converts it into machine language which is a binary representation of the assembly using the instruction formats described in the ISA. More advanced discussion:

The ~~Assembler~~ linker patches together separately assembled object files into one executable and the loader loads the executable into memory so the processor can begin execution.

2. a) die size is $18.9\text{mm} \times 13.6\text{mm}$
wafer diameter is 300mm

The following equation was provided in the lecture slides:

$$\begin{aligned} \text{Dies per wafer} &= \frac{\pi \times (\text{Wafer diameter}/2)^2}{\text{Die area}} - \frac{\pi \times \text{Wafer diameter}}{\sqrt{2 \times \text{Die area}}} \\ &= \frac{\pi \times (300\text{mm}/2)^2}{18.9\text{mm} \times 13.6\text{mm}} - \frac{\pi \times 300\text{mm}}{\sqrt{2 \times 18.9\text{mm} \times 13.6\text{mm}}} \end{aligned}$$

$$= 275.00 - 41.568 = 233.43 \text{ dies}$$

so 233 dies

2 b)

$$\begin{aligned} \text{die yield} &= \text{wafer yield} \times \frac{1}{\left(1 + \frac{\text{defects}}{\text{area}} \times \text{die area}\right)^N} \\ &= 0.99 \times \frac{1}{\left(1 + \frac{0.020 \text{ defects}}{\text{cm}^2} \times \frac{1 \text{ cm}^2}{100 \text{ mm}^2} \times 18.9 \text{ mm} \times 13.6 \text{ mm}\right)^7} \\ &= 0.697 \approx 69.7\% \end{aligned}$$

c)

$$\begin{aligned} \text{cost per die} &= \frac{\text{cost per wafer}}{\text{dies per wafer} \times \text{die yield}} \\ &= \frac{\$6500}{233 \text{ dies} \times 69.7\%} = \$40.02 \text{ per die} \end{aligned}$$

d)

This problem is really just

$$\begin{array}{c} \text{profit} = 1.5x - 1.0x = 0.5x + 6,500 = \$3,250 \text{ per wafer} \\ \uparrow \qquad \qquad \uparrow \\ \text{price} \qquad \text{cost} \end{array}$$

e)

$$\begin{array}{c} 233 \text{ dies per wafer} \times (1 - 0.697) \times (0.65) \times (0.6) \times \$40.02 = \$1,652.80 \\ \uparrow \qquad \qquad \uparrow \qquad \qquad \uparrow \qquad \qquad \uparrow \\ \text{failed dies} \quad \% \text{ of failed dies repurposed} \quad \% \text{ of full price repurposed dies can be sold for} \end{array}$$

$\$40.02 \times 1.5 = \$1,652.80$

$$\text{total gains per wafer} = \$3,250.00 + 1,652.80 = \$4,902.80$$

$$\text{left over after R\&D} = 0.85 \times \$4,902.80 = \$4,169.40$$

3

a) P₃ has the highest clock rate @ 3.26Hz
clock period of each:

$$P_{1 \text{ period}}^{\text{clock}} = \frac{1}{1.86\text{Hz}} = 0.535\text{ns} = 535\text{ps}$$

$$P_{2 \text{ period}}^{\text{clock}} = \frac{1}{2.46\text{Hz}} = 0.4167\text{ns} = ~~416.7\text{ps}~~ 416.7\text{ps}$$

$$P_{3 \text{ period}}^{\text{clock}} = \frac{1}{3.26\text{Hz}} = 0.3125\text{ns} = 312.5\text{ps}$$

b)

$$\text{MIPS} = \frac{\text{clock rate}}{\text{CPI}} \cdot \frac{1}{10^6}$$

$$\text{MIPS}_{P_1} = \frac{1.86\text{Hz}}{0.8\text{cc/instr}} \cdot 10^{-6} = 2,250\text{MIPS}$$

$$\text{MIPS}_{P_2} = \frac{2.46\text{Hz}}{1.2\text{cc/instr}} \cdot 10^{-6} = 2,000\text{MIPS}$$

$$\text{MIPS}_{P_3} = \frac{3.26\text{Hz}}{1.6\text{cc/instr}} \cdot 10^{-6} = 2,000\text{MIPS}$$

P₁ has the highest MIPS

c) $t_{\text{execution}} = \frac{\text{dynamic instr. count}}{\text{instr./sec}}$

$$t_{\text{ex } P_1} = \frac{22.0 \text{ million instr}}{2,250 \text{ MIPS}} = \underline{9.78 \text{ ms}}$$

$$t_{\text{ex } P_2} = \frac{19 \text{ million instr}}{2,000 \text{ MIPS}} = \underline{9.5 \text{ ms}} \quad \underline{P_2 \text{ is faster!}}$$

$$t_{\text{ex } P_3} = \frac{22.0 \text{ million instr}}{2,000 \text{ MIPS}} = \underline{11.0 \text{ ms}}$$

3 d) the reference machine takes 35ms

$$\text{speedup}_1 = \frac{35\text{ms}}{9.78\text{ms}} = \underline{3.58x}$$

$$\text{speedup}_2 = \frac{35\text{ms}}{9.5\text{ms}} = \underline{3.68x}$$

$$\text{speedup}_3 = \frac{35\text{ms}}{11.0\text{ms}} = \underline{3.18x}$$

4. a) This system is I/O bound because it spends more than 50% of its time waiting on I/O

b)
$$t_{\text{improved}} = \frac{t_{\text{affected}}}{\text{speedup}_{\text{affected}}} + t_{\text{unaffected}}$$

let's assume the program originally took a second

$$t_{\text{improved}} = \frac{(0.4)(1s)}{20} + (0.6)(1s) \\ = 0.62s$$

$$\text{speedup} = \frac{1s}{0.62s} = \underline{1.61x}$$

c)
$$t_{\text{improved}} = \frac{0.6s}{3} + 0.4s \\ = 0.60s$$

$$\text{speedup} = \frac{1s}{0.6s} = \underline{1.67x}$$

d) you can speedup disk access by using a parallel disk array, using a disk with a higher RPM, or by replacing disk with flash.

$$5. \quad a) \quad t_{\text{new}} = \frac{800s}{4} + 200s = \underline{400s}$$

$$\text{speedup} = \frac{1000s}{400s} = \underline{2.5x}$$

$$b) \quad t_{\text{new}} = \frac{800s}{8} + 200s = \underline{300s}$$

$$\text{speedup} = \frac{1000s}{300s} = \underline{3.33x}$$

$$c) \quad t_{\text{new}} = \frac{800s}{20} + 200s = \underline{200s}$$

$$\text{speedup} = \frac{1000s}{200s} = \underline{5x}$$

$$d) \quad t_{\text{new}} = \frac{900s}{20} + 100s = \underline{100s}$$

$$\text{speedup} = \frac{1000s}{100s} = \underline{10x}$$

6.

$$a) \quad t_{\text{new}} = \frac{2s / (0.95)}{100} + 2s(1-0.95)$$

$$= 0 + 0.1s$$

$$\text{speedup} = \frac{2s}{0.1s} = \underline{20x}$$

$$b) \quad t_{\text{new}} = \frac{2s(0.95)}{50} + 2s(1-0.95)$$

$$= \underline{0.138s}$$

$$\text{speedup} = \frac{2s}{0.138s} = \underline{14.493x}$$

6. c) Energy = power \times time

the time will be reduced by 14.493x

while the power will be doubled

$$\frac{\text{Energy old}}{\text{Energy new}} = \frac{14.493}{2} = \underline{7.25x}$$

d) see attached graph and GNU plot source file for generating the graph

7. a)

$$P_{\text{total}} = P_{\text{dynamic}} + P_{\text{static}}$$
$$= 1.2 \text{ W}$$

$$\text{Energy per operation} = P_{\text{total}} \times t_{\text{operation}}$$
$$= 1.2 \text{ W} \times 1 \text{ s}_{\text{op}} = 1.2 \text{ J}_{\text{op}}$$

b) In this case dynamic power will change as well as the operation time

$$P_{\text{dynamic}} \propto C V^2 f \text{ thus } P_{\text{new dynamic}} = 1.0 \text{ W} \cdot \frac{250 \text{ MHz}}{500 \text{ MHz}}$$

Note that static power does not depend on f and thus does not change $= 0.5 \text{ W}$

$$t_{\text{new}} = 1 \text{ s} \cdot \frac{500 \text{ MHz}}{250 \text{ MHz}} = 2 \text{ s}$$

$$P_{\text{total}} = 0.5 \text{ W} + 0.2 \text{ W} = \underline{0.7 \text{ W}}$$

$$\text{Energy}_{\text{op}} = (0.5 \text{ W} + 0.2 \text{ W})(2 \text{ s}) = \underline{1.4 \text{ J/s}}$$

notice the energy per op went up due to longer execution time burning more static energy!

7. c) In this case, we also reduce the voltage from 1.2V to 0.9V, taking advantage of the V^2 relationship to dynamic power

$$P_{\text{dynamic}} = 1.0W \times \frac{250\text{MHz}}{500\text{MHz}} \times \left(\frac{0.9V}{1.2V}\right)^2 = 0.281W$$

$$P_{\text{static}} = V \cdot I_{\text{leakage}} = 0.2W \times \frac{0.9V}{1.2V} = 0.15W$$

$$P_{\text{total}} = 0.281W + 0.15W = \underline{0.431W}$$

$$\text{Energy per op} = 0.431W \times 25 = \underline{0.863J/op}$$

Now we see an energy savings! but only when we drop voltage as well as frequency!

- d) This problem is the ~~the~~ reverse of the previous.

$$P_{\text{dynamic}} = 1.0W \times \frac{650\text{MHz}}{500\text{MHz}} \times \left(\frac{1.4V}{1.2V}\right)^2 = 1.7694W$$

$$P_{\text{static}} = 0.2W \times \frac{1.4V}{1.2V} = 0.2333W$$

$$t_{\text{new}} = 15 \times \frac{500\text{MHz}}{650\text{MHz}} = 0.7695$$

$$P_{\text{total}} = 1.7694W + 0.2333W = \underline{2.003W}$$

$$\text{Energy per op} = 2.00W \times 0.7695 = \underline{1.54J/op}$$

- e) Power gating is one technique that removes power from unused logic. For example, an unused floating point unit or even an entire core might be powered down. This saves dynamic and static energy!

Clock gating disables the clock to unused logic. This only saves dynamic energy but keeps static elements powered

7. e) continued...

such as memory and registers. This allows a given logic block to retain state while saving power.

Clock gating does not save as much energy as power gating but allows logic to come alive more quickly.

Other techniques include ~~hardware~~ hardware acceleration where commonly used operations are mapped into more energy efficient hardware.

A modern system on a chip (SOC) will use all techniques discussed to save energy!

f) The previous problems changed V and f but left C alone. This problem changes all three!

$$P_{\text{dynamic}} = \frac{1.0W}{1.5W} \left(\frac{130nm}{90nm} \right) \left(\frac{1.5V}{1.2V} \right)^2 \left(\frac{333MHz}{500MHz} \right) = \underline{1.5W}$$

↑ ↑ ↑
capacitance increase voltage increase frequency goes down!

$$P_{\text{static}} = \left(\frac{1.5V}{1.2V} \right) 0.2W = 0.25W$$

$$P_{\text{total}} = 1.5W + 0.25W = \underline{1.75W}$$

$$t_{\text{new}} = \cancel{1.5} \times \left(\frac{333MHz}{500MHz} \right) = \leftarrow \text{time increases!}$$
$$= 1.5 \times \frac{500MHz}{333MHz} = 1.5015s$$

$$\text{Energy per op} = \cancel{1.5W} \times \cancel{1.5015s}$$
$$1.75W \times 1.5015s = \underline{2.63J/op}$$
$$\frac{2.63J/op}{1.2J/op} = \underline{2.19x \text{ higher!}}$$

8. (10 points) The execution times of a few SPEC2006Cint benchmarks on three different machines are shown in the table below.

Benchmark	Ultra 5 Time (sec)	Machine A Time (sec)	Machine A SPEC ratio	Machine B Time (sec)	Machine B SPEC ratio
perlbench	9,770	454	21.52	253	38.62
bzip2	9,650	520	18.56	438	22.03
gcc	8,050	461	17.46	234	34.40
mcf	9,120	268	34.03	150	60.80
gobmk	10,490	429	24.45	383	27.39
hammer	9,330	197	47.36	135	69.11
sjeng	12,100	529	22.87	362	33.43
libquantum	20,720	91	227.69	7	2960.00
h264ref	22,130	599	36.95	427	51.83
omnetpp	6,250	305	20.49	163	38.34
astar	7,020	327	21.47	234	30.00
xalancbmk	6,900	253	27.27	117	58.97
Geometric mean			30.45		56.98

- (a) Using the Ultra 5 as a baseline, fill in the SPEC ratio columns above for Machines A and B.
 (b) Compute the geometric mean of the SPEC ratios for both machines. *see table above*
 (c) Given the benchmarks above, which machine (A or B) is faster? By how much?
 (d) libquantum has fallen victim to substantial compiler optimizations as evidenced by such extraordinary speedups over the reference machine. Consequently, it has been removed from SPEC2017. Recalculate the geometric mean of the normalized performance without libquantum. Compare machine A to machine B excluding libquantum.

c) Machine B is faster because it has the higher spec ratio. Specifically Machine B is faster than Machine A by

$$\frac{56.98}{30.45} = 1.87x$$

d)

Machine A
Geometric mean: 25.36

Machine B
Geometric mean: 39.79

Machine B is still faster but by a lesser amount:

$$\frac{39.79}{25.36} = 1.57x$$

9. a) $CPI = \text{clock cycles per instr}$

Program A instr count = $8 \cdot 10^8$ instr
 execution time = 1.1s

Program B instr count = $1.2 \cdot 10^9$ instr
 execution time = 1.6s

both run on processor w/ $\frac{1}{\text{ns}} = 1 \text{ GHz}$ clock rate
using units to derive equation:

$$CPI = \text{cc/instr} = \frac{\text{cc}}{\text{s}} \times \text{s} \times \frac{1}{\text{instr}}$$

$$CPI_A = 10^9 \text{ cc/s} \times 1.1 \text{ s} \times \frac{1}{8 \cdot 10^8 \text{ instr}} = \frac{1.1}{8} \cdot 10 \text{ cc/instr}$$
$$= \underline{1.375 \text{ cc/instr}}$$

$$CPI_B = 10^9 \text{ cc/s} \times 1.6 \text{ s} \times \frac{1}{1.2 \cdot 10^9 \text{ instr}} = \frac{1.6}{1.2} \text{ cc/instr}$$
$$= \underline{1.333 \text{ cc/instr}}$$

b) $t_{\text{execute}} = \frac{CPI \times \text{instr. count}}{\text{clock rate}}$

or
 $\text{clock rate} = \frac{CPI \times \text{instr count}}{t_{\text{execute}}}$

$$\frac{\text{clock rate}_A}{\text{clock rate}_B} = \frac{CPI_A \times \text{instr count}_A}{t_{\text{execute}_A}} \cdot \frac{t_{\text{execute}_B}}{CPI_B \times \text{instr count}_B}$$

$$t_{\text{ex}_A} = t_{\text{ex}_B} \text{ so...} = \frac{CPI_A \times \text{instr count}_A}{CPI_B \times \text{instr count}_B}$$

9. b) continued...

$$\frac{\text{clock rate}_A}{\text{clock rate}_B} = \frac{1.375 \text{ cc/instr} \times 8 \cdot 10^8 \text{ instr}}{1.333 \text{ cc/instr} \times 1.2 \cdot 10^9 \text{ instr}}$$
$$= \underline{0.688}$$

c) original processor ran a 164x

$$t_{ex_c} = \frac{CPI_c \times \text{instr. count}_c}{\text{clock rate}}$$

$$\text{speedup}_{cA} = \frac{t_{exA}}{t_{ex_c}} = \frac{CPI_A \times \text{instr. count}_A}{\text{clock rate}} \times \frac{\text{clock rate}}{CPI_c \times \text{instr. count}_c}$$
$$= \frac{1.375 \text{ cc/instr} \times 8 \cdot 10^8 \text{ instr}}{1.3 \text{ cc/instr} \cdot 7.5 \cdot 10^8 \text{ instr}} = \underline{1.13x}$$

$$\text{speedup}_{cB} = \frac{1.333 \text{ cc/instr} \cdot 1.2 \cdot 10^9 \text{ instr}}{1.3 \text{ cc/instr} \cdot 7.5 \cdot 10^8 \text{ instr}} = \underline{1.64x}$$